

Reactor

From complexity to clarity: demystifying agent development

Cihan Cinar Software & Al Architect





Cihan Cinar Software & Al Architect Givaudan

Linkedin: @cihancinar Github: @cihancinar X: @Cihan_Cinar_

https://blog.cihan-cinar.com/



Agentic AI Terms You Should Know!



What is Generative AI

Generative AI is a branch of artificial intelligence focused on creating new content. Unlike traditional AI, which might classify or predict based on data, generative AI can produce original outputs such as images, text, music and more based on the vast amounts of data that the generative model was trained on.

Microsoft



What are large language models (LLMs)

Large language models (LLMs) are advanced AI systems that understand and generate natural language, or human-like text, using the data they've been trained on through machine learning techniques.



Microsoft



Chat with LLM

(i)



Limited to the LLMs knowledge and user input

Hicrosoft



Elements of a Good Prompt

Act as a helpful tutor who breaks down complex subjects into easy explanations.

I want you to explain the process of photosynthesis to a

14 year old student, to assist with biology exam preparations.

Your answer should be 300 words, written in a tone that's friendly and educational.

Persona: Ask the tool to take a role

Objective: What do you want the AI to do

Audience: Specify who it's for Context: What does the tool need to know

Boundaries: Set your own direction & limitation

Tip 1 Give Clear Instructions

Use commands that instruct the Al tool on what you want to generate, such as 'explain', 'translate', 'summarize' or 'compare'.

Tip 2 Provide Context

Adding context and background information can help the tool to understand the task better. For example, mention the project type such as 'short story', 'report' or 'outline

Tip 3 Iterate & Experiment

Try different instructions and techniques if you don't get the results you want. Prompting can be like an experiment that may require several rounds of iterations!

Chat Completion

- conversation-in and message-out
- LLM based on token



- Context window
 - 4k input & 4k output tokens for common models
 - 128k+ input & 16k+ output tokens for newer models



Chat Completion

Request

Response

```
POST
https://{endpoint}/openai/deployments/{deploymen
                                                                           "choices": [
t-id}/chat/completions?api-version=2024-10-21
                                                                               "finish_reason": "stop",
                                                                               "index": 0,
                                                                               "message": {
    "model": "gpt-40",
                                                                                 "content": "Aye! Provide fresh food n' water, a roomy cage,
    "messages": [
                                                                       toys fer enrichment, regular vet check-ups, and plenty o' social
                                                                       interaction. Arr!",
        "role": "system",
                                                                                 "role": "assistant"
        "content": "you are a helpful assistant that talks like
a pirate. Generate short answer only."
      },
                                                                            "created": 1740803441,
        "role": "user",
                                                                           "id": "chatcmpl-B68mnFNQafMCjwosWz8zPS1ZtcRVO",
        "content": "can you tell me how to care for a parrot?"
                                                                           "model": "gpt-4o-2024-05-13",
                                                                            "object": "chat.completion",
                                                                            "system fingerprint": "fp_65792305e4",
    ر [
    "temperature": 0.7,
                                                                           "usage": {
    "stream": false,
                                                                              "completion tokens": 32,
    "max tokens": 2000
                                                                              "prompt tokens": 39,
                                                                             "total tokens": 71
```

LLM Chat with RAG

(i)



Relevant answers, but limited to data sources

Hicrosoft



How RAG Works (Retrieval-Augmented Generation)

- Document ingestion
- Transform, the user query into a vector
- Search in vector DB for relevant documents
- Add found documents to the context

Microsoft





• Transform, the user query into a vector



Microsoft



• Search in vector DB for relevant documents



• Add found documents to the context



Al Agent

An Agent is a semi-autonomous software that leverages large language models to operate independently over extended periods, using various tools to accomplish complex task.

It can be defined as a prescriptive implementation that follows predefined workflows.

All of these variations are considered agentic systems.

Agentic systems



Reactor

Building block

The basic building block of agentic systems is an LLM enhanced with augmentations such as retrieval, tools, and memory.



e.g. Get weather from an API tools.

- Microsoft



Chain Workflow

The Chain Workflow pattern exemplifies the principle of breaking down complex tasks into simpler, more manageable steps.



e.g. Generate a marketing copy, then translating it into a different language

- Microsoft



Parallelization Workflow

e.g.

LLMs can work simultaneously on tasks and have their outputs aggregated programmatically.



Review code with different profile: Security, Performance, Quality



Routing Workflow

The Routing pattern implements intelligent task distribution, enabling specialized handling for different types of input.



e.g.

Switch between LLM like for coding, reasoning, image generation, lower cost LLM

- Microsoft

Orchestrator-Workers Workflow

In the orchestrator-workers workflow, a central LLM dynamically breaks down tasks, delegates them to worker LLMs, and synthesizes their results.

e.g. Search data from database, web, LLM and merge them

- Microsoft

Evaluator-Optimizer

The Evaluator-Optimizer pattern implements a dual-LLM process where one model generates responses while another provides evaluation and feedback in an iterative loop, similar to a human writer's refinement process.

Microsoft

LLM Chat with Agent

Knowledge sources

i Agent perform complex tasks

Hicrosoft

LLM Chat with Multi Agent

Knowledge sources

(i) Specific assigned task agent

Hicrosoft

Al Agent Service in Action

Step 1 Create an Agent	Agent Microsoft Sales Agent	Instruction You are an advanced sales analyses agent for Microsoft,	Thread Sales analysis	Run 1
Step 2 Create a Thread		specializing in assisting users with sales data inquiries	Wer's message Tell me the total sales by region	Function Calling Tool Query SQLite DB
Step 3 Run the Agent		Model		
Step 4 Check the Run status		Your data (optional)	Here is the sales: Europe: \$1923 America: \$1776 Run 2	
Step 5 Display the Agent's Response		Azure Al Search Files (local or Azure Blob)		Run 2
		Tools (optional) File Search Code Interpreter Function Calling Bing Search Microsoft SharePoint Microsoft Fabric Azure Logic Apps Azure Functions OpenAl 3.0 specified tools	Show a pie chart	Code Interpreter Tool Create a pie chart Create message
			Agent's message	

SDK

Language	SDK
Java	Azure SDK for Java, Semantic Kernel, Spring AI, LangChain4j
Python	Azure SDK for Python, Semantic Kernel, AutoGen, LangChain
C#	Azure SDK for .NET, Semantic Kernel
JS	Azure SDK for Js, LangChainJs, Vercel AI SDK

Microsoft Reactor

Demo

Agent

Is RAG an Agent?

Resources

- <u>https://github.com/cihancinar/nextjs-azure-ai-starter</u>
- <u>https://aka.ms/Apr7AzureAlServicesOpenAl1</u>
- <u>https://aka.ms/Apr7AlFoundryConceptsRAG1</u>
- <u>https://github.com/Azure-Samples/</u>
- <u>https://github.com/Azure/GPT-RAG</u>
- <u>https://www.anthropic.com/engineering/building-effective-agents</u>

Hicrosoft

Reactor

Thank You

Linkedin: @cihancinar Github: @cihancinar X: @Cihan_Cinar_

https://blog.cihan-cinar.com/

Cihan Cinar Software & Al Architect

